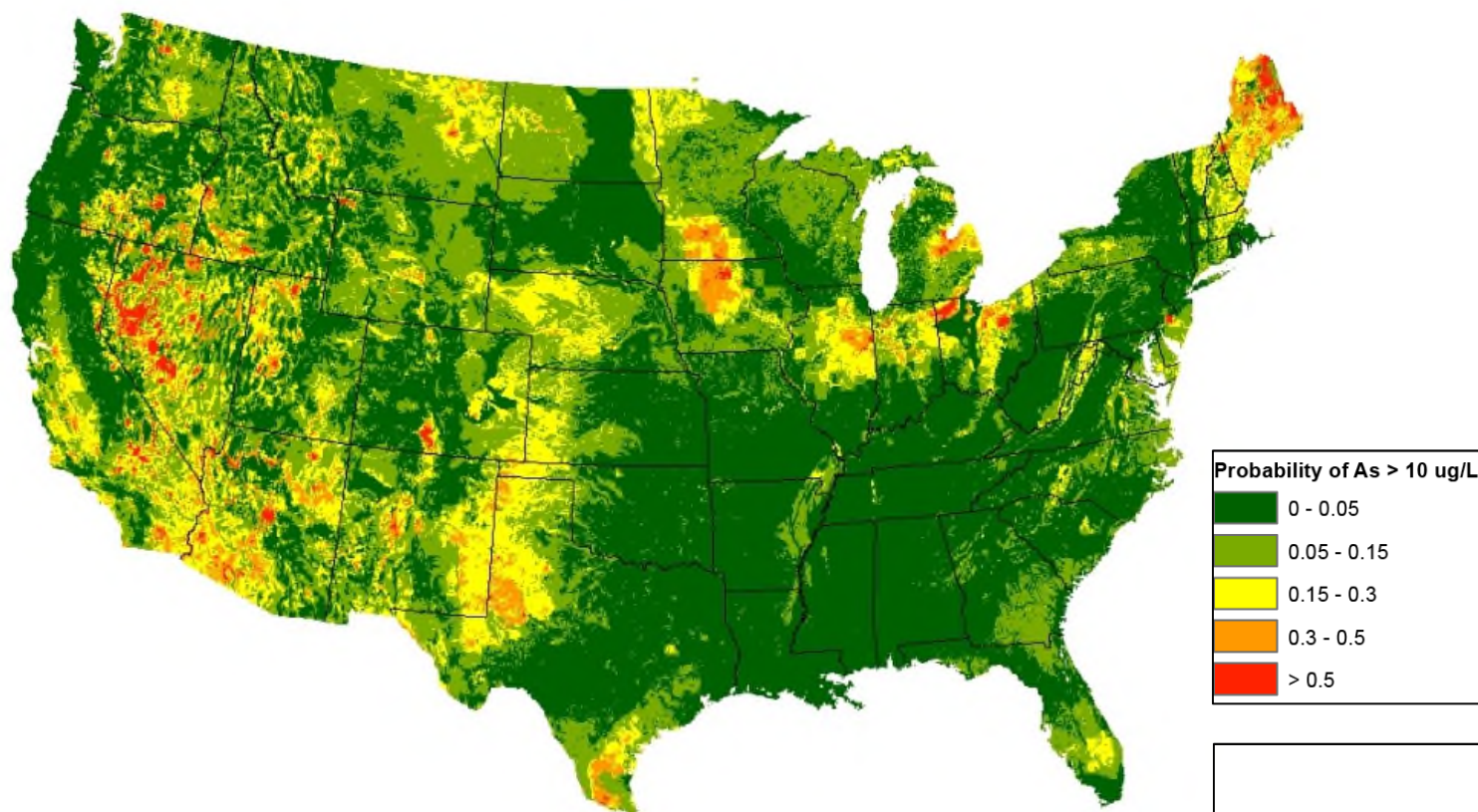# A comparison of statistical modeling techniques to predict arsenic in domestic wells in the CONUS

Melissa Lombard, Bernard T. Nolan, Mathew Gribble, Maria Argos, Joseph Ayotte

## Estimating the High-Arsenic Domestic-Well Population in the Conterminous United States

Joseph D. Ayotte,[*,†] Laura Medalie,[‡] Sharon L. Qi,[§] Lorraine C. Backer,[∥] and Bernard T. Nolan[⊥]

Probability of As > 10 ug/L
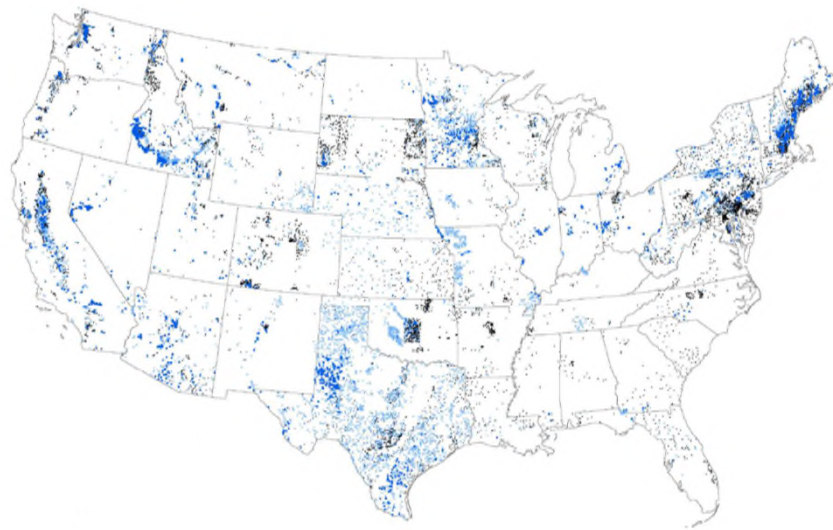- 0 - 0.05
- 0.05 - 0.15
- 0.15 - 0.3
- 0.3 - 0.5
- > 0.5

- Logistic Regression model

- Model response term is the probability of arsenic > 10μg/L

$$P(y|x) = \frac{e^{\beta_0+\beta_1 x_1 +\cdots+\beta_n x_n}}{1 + e^{\beta_0+\beta_1 x_1 +\cdots+\beta_n x_n}}$$

USGS
science for a changing world

# Existing logistic regression model

Arsenic concentrations
from 20,450 domestic wells



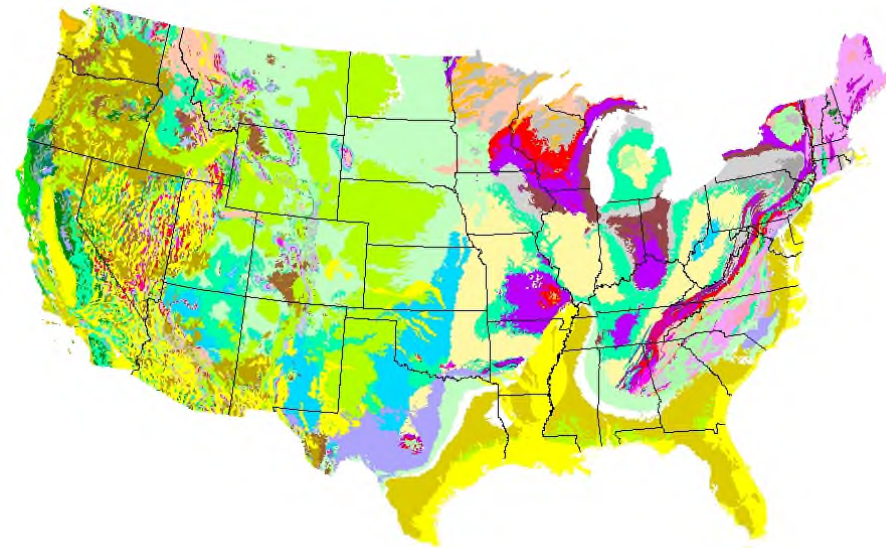Black = As < 1 $\mu$g/L
Light blue = 1 $\geq$ As < 10 $\mu$g/L
Dark blue = As > 10 $\mu$g/L

*From Ayotte et al. 2017, ES&T*
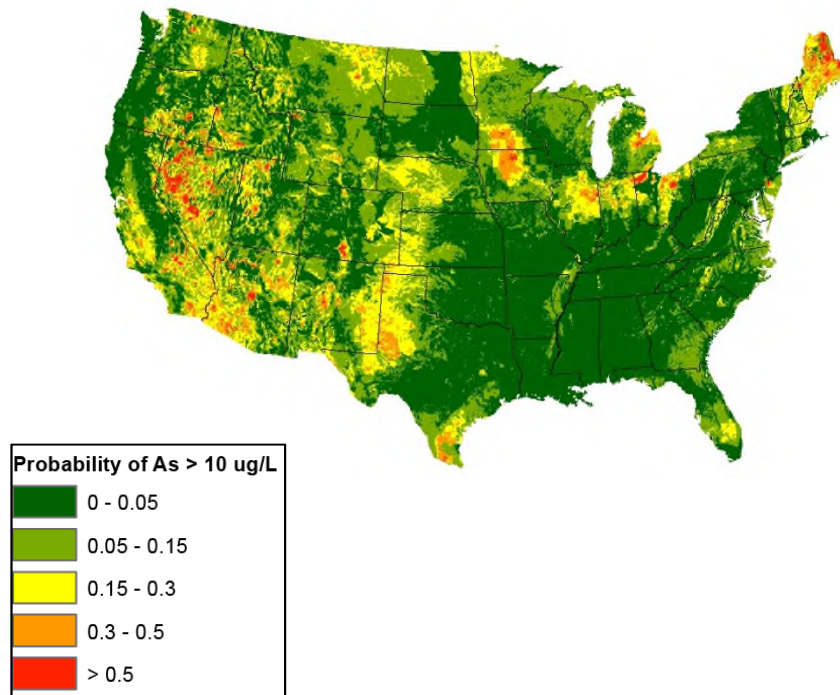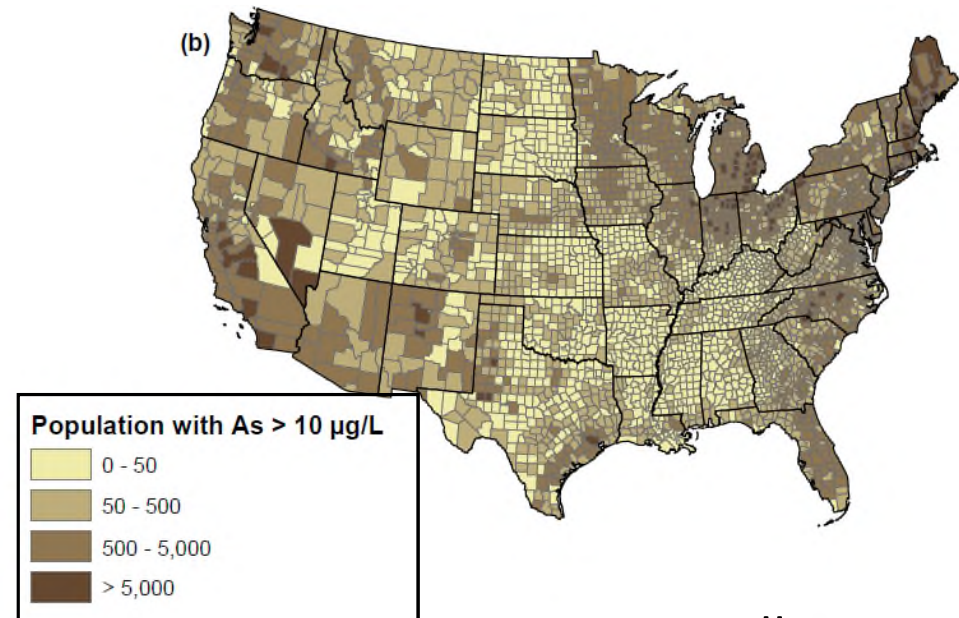
42 predictor variables



Generalized geology of the U.S.

*After King and Beikman, 1974*

# Existing arsenic model and exposure estimate

Arsenic model output

Arsenic exposure estimate



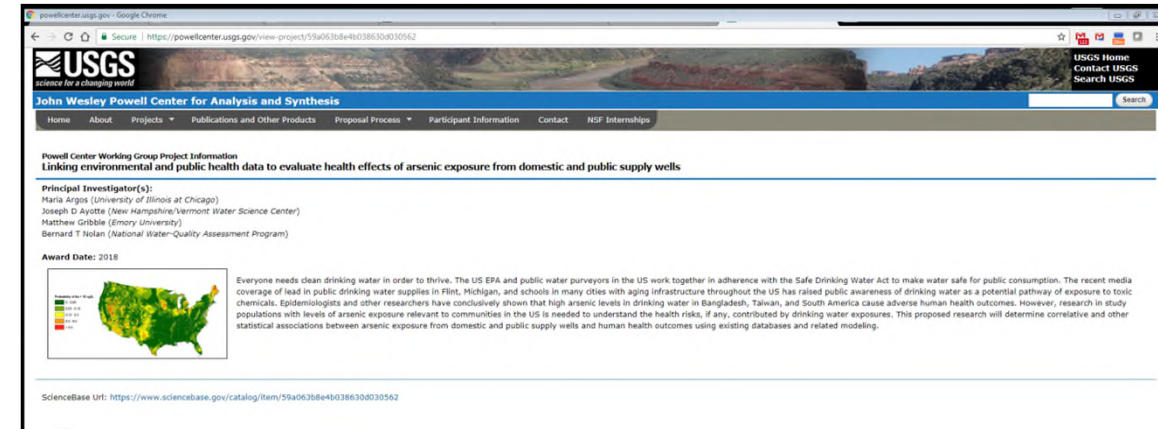Probability of As > 10 ug/L
- 0 - 0.05
- 0.05 - 0.15
- 0.15 - 0.3
- 0.3 - 0.5
- > 0.5

Population with As > 10 µg/L
- 0 - 50
- 50 - 500
- 500 - 5,000
- > 5,000

2.1 million people

*From Ayotte et al. 2017, ES&T*

**≋USGS**
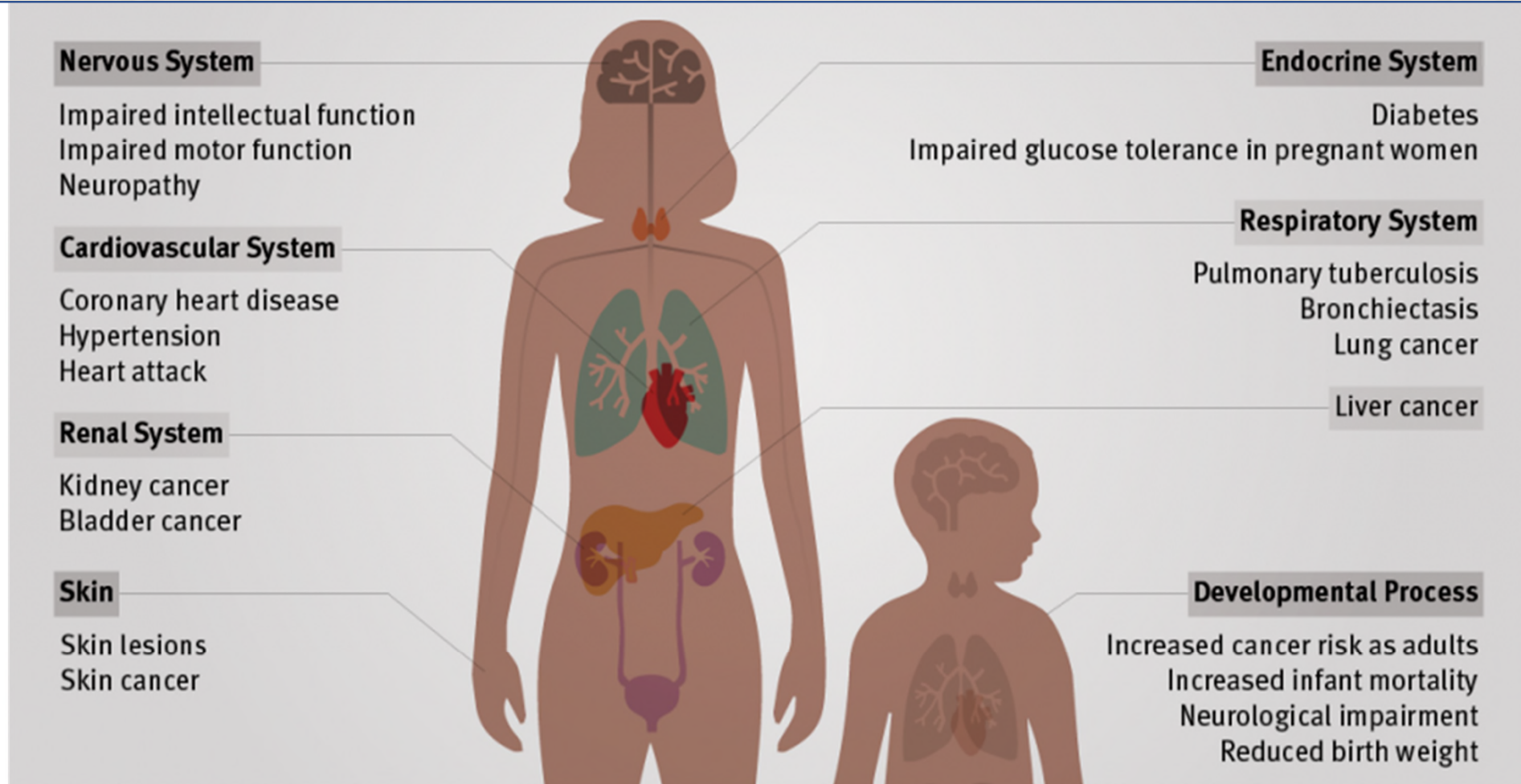*science for a changing world*

# USGS Powell Center study

Linking environmental and public health data to evaluate health effects of arsenic exposure from domestic supply wells

- Make a new national arsenic model
  - Use machine learning methods
  - Update model variables



- Results will be used in epidemiology models to evaluate relationships between human health outcomes and arsenic in domestic wells.
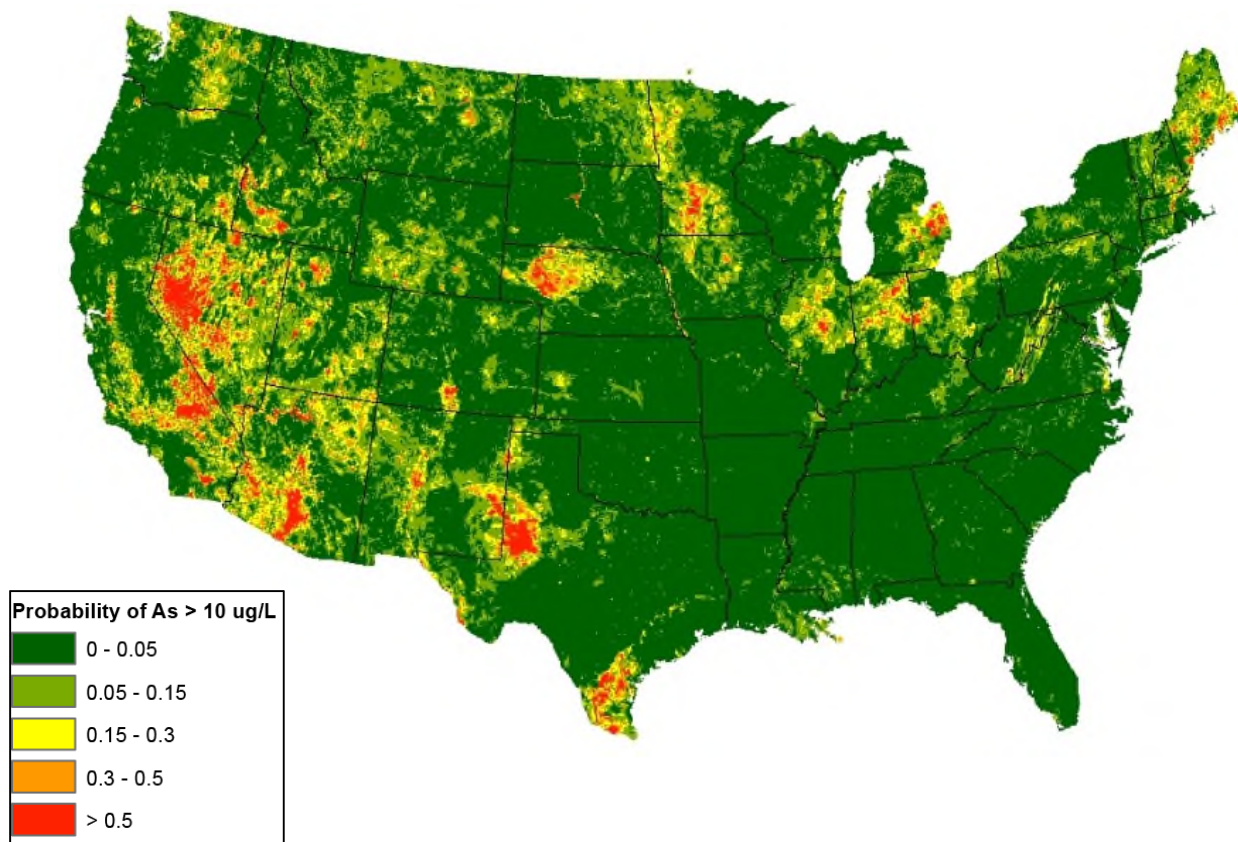
≋USGS
science for a changing world

# Arsenic's effects on the human body



**Nervous System**

Impaired intellectual function
Impaired motor function
Neuropathy

**Cardiovascular System**

Coronary heart disease
Hypertension
Heart attack

**Renal System**

Kidney cancer
Bladder cancer

**Skin**

Skin lesions
Skin cancer

**Endocrine System**

Diabetes
Impaired glucose tolerance in pregnant women

**Respiratory System**

Pulmonary tuberculosis
Bronchiectasis
Lung cancer

Liver cancer

**Developmental Process**

Increased cancer risk as adults
Increased infant mortality
Neurological impairment
Reduced birth weight

https://www.hrw.org/report/2016/04/06/nepotism-and-neglect/failing-response-arsenic-drinking-water-bangladeshs-rural#page

USGS
*science for a changing world*

# Boosted Regression Tree Model



**Probability of As > 10 ug/L**
- 0 - 0.05
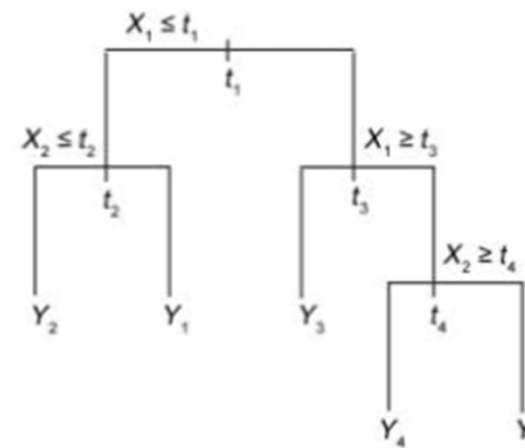- 0.05 - 0.15
- 0.15 - 0.3
- 0.3 - 0.5
- > 0.5

BRT Model
Number of trees = 4500
Interaction depth = 14
Learning rate = 0.01



*From Elith et al., 2008*

# LR and BRT Model Comparison

LR Model

Probability of As > 10 ug/L
- 0 - 0.05
- 0.05 - 0.15
- 0.15 - 0.3
- 0.3 - 0.5
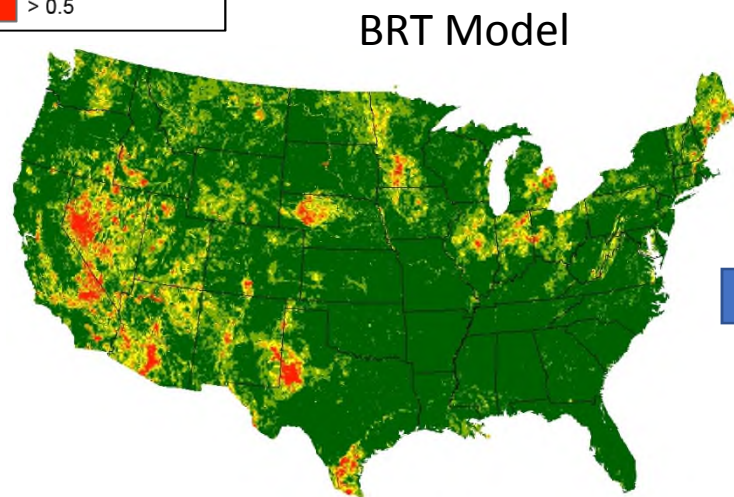- > 0.5

BRT Model

- Sharper delineation between high and low arsenic probabilities with BRT model

BRT minus LR probability

Change in Probability
- -1 - -0.4
- -0.4 - -0.2
- -0.2 - 0
- 0 − 0.4
- 0.4 - 1

≋USGS
science for a changing world

# Model Comparison Predictive Performance

|  | Training Data | | | Hold-out Data | | |
|---|---|---|---|---|---|---|
|  | Accuracy | Sensitivity | Specificity | Accuracy | Sensitivity | Specificity |
| **LR** | 89.9% | 12.7% | 99.3% | 90.1% | 13.9% | 99.0% |
| **BRT** | 95.6% | 64.5% | 99.3% | 92.1% | 43.7% | 97.8% |

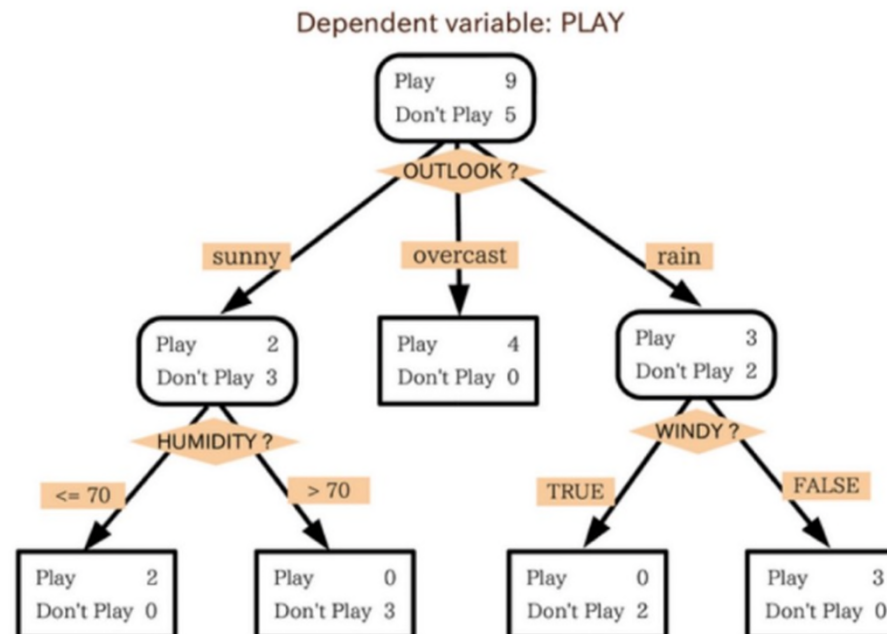Sensitivity = events  (As>10μg/)          Specificity = non-events (As<10μg/L)

# Random Forest Classification Model

- Ensemble tree based machine learning method

- Model response term is a classification (category)
  - Arsenic ≤ 10 μg/L
  - Arsenic > 10 μg/L
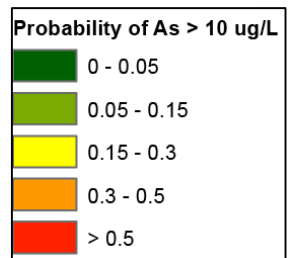
RFC Model
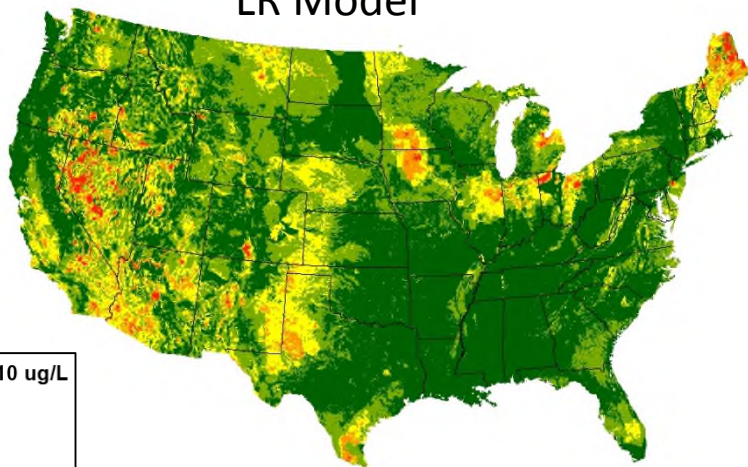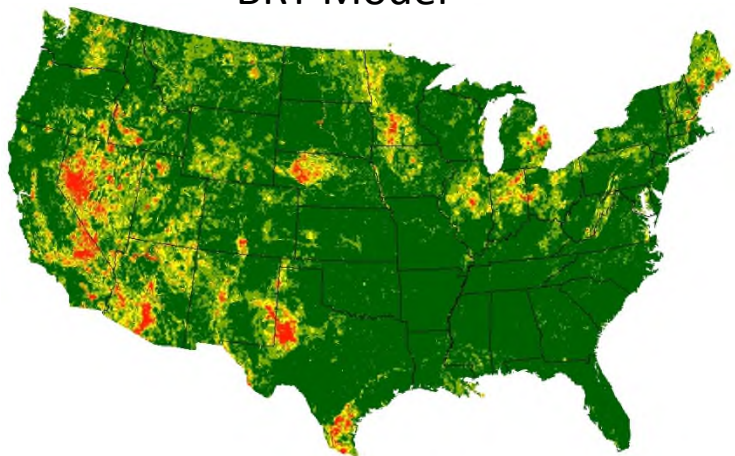Number of trees = 500
mtry = 38

Dependent variable: PLAY



From medium.com

# LR, BRT, & RFC Model Comparison



LR Model

BRT Model

RFC Model

Probability of As > 10 ug/L
- 0 - 0.05
- 0.05 - 0.15
- 0.15 - 0.3
- 0.3 - 0.5
- > 0.5

Legend
- 1 Less than 10 ug/L
- 2 Greater than 10 ug/L

# Model Comparison Predictive Performance

| | Training Data | | | Hold-out Data | | |
|---|---|---|---|---|---|---|
| | Accuracy | Sensitivity | Specificity | Accuracy | Sensitivity | Specificity |
| **LR** | 89.9% | 12.7% | 99.3% | 90.1% | 13.9% | 99.0% |
| **BRT** | 95.6% | 64.5% | 99.3% | 92.1% | 43.7% | 97.8% |
| **RFC 2C** | 99.9% | 100% | 99.3% | 91.0% | 38.4% | 97.1% |

Sensitivity = events (As>10μg/L )
Classification As>10μg/L

Specificity = non-events (As<10μg/L)
Classification As<10μg/L

USGS
science for a changing world

# Next steps

- Developing a RFC model with 4 concentration categories
  - ≤ 5 µg/L
  - 5 – 10 µg/L
  - 10 – 50 µg/L
  - > 50 µg/L

- Arsenic model results will be used in epidemiology models
  - Low birth weight
  - Cancers
  - Diabetes
  - Cardiovascular disease



Cancer Mortality Rates by State Economic Area (Age-adjusted 1970 US Population)
Bladder: White Males, 1970-94

US = 6.56/100,000

7.92-10.84 (highest 10%)
7.25- 7.91
6.81- 7.24
6.45- 6.80
6.14- 6.44
5.82- 6.13
5.41- 5.81
5.00- 5.40
4.68- 4.99
3.16- 4.67 (lowest 10%)

NATIONAL CANCER INSTITUTE®

≋USGS
science for a changing world

CONTACT INFO:

Melissa Lombard, Ph.D.
US Geological Survey
331 Commerce Way
Pembroke, NH 03275
mlombard@usgs.gov
603-226-7816

# Questions ?

![USGS logo - science for a changing world]